

Column Subset Selection, Matrix Factorization, and Eigenvalue Optimization

Joel A. Tropp*

26 June 2008. Revised 2 October 2008.

Abstract

Given a fixed matrix, the problem of column subset selection requests a column submatrix that has favorable spectral properties. Most research from the algorithms and numerical linear algebra communities focuses on a variant called rank-revealing QR, which seeks a well-conditioned collection of columns that spans the (numerical) range of the matrix. The functional analysis literature contains another strand of work on column selection whose algorithmic implications have not been explored. In particular, a celebrated result of Bourgain and Tzafriri demonstrates that each matrix with normalized columns contains a large column submatrix that is exceptionally well conditioned. Unfortunately, standard proofs of this result cannot be regarded as algorithmic. This paper presents a randomized, polynomial-time algorithm that produces the submatrix promised by Bourgain and Tzafriri. The method involves random sampling of columns, followed by a matrix factorization that exposes the well-conditioned subset of columns. This factorization, which is due to Grothendieck, is regarded as a central tool in modern functional analysis. The primary novelty in this work is an algorithm, based on eigenvalue minimization, for constructing the Grothendieck factorization. These ideas also result in an approximation algorithm for the $(\infty, 1)$ norm of a matrix, which is generally NP-hard to compute exactly. As an added bonus, this work reveals a surprising connection between matrix factorization and the famous MAXCUT semidefinite program.

1 Introduction.

Column subset selection refers to the challenge of extracting from a matrix a column submatrix that has some distinguished property. These properties commonly involve conditions on the spectrum of the submatrix. The most familiar example is probably rank-revealing QR, which seeks a well-conditioned collection of columns that spans the (numerical)

range of the matrix [GE96].

The literature on geometric functional analysis contains several fundamental theorems on column subset selection that have not been discussed by the algorithms community or the numerical linear algebra community. These results are phrased in terms of the *stable rank* of a matrix:

$$\text{st.rank}(\mathbf{A}) = \frac{\|\mathbf{A}\|_{\text{F}}^2}{\|\mathbf{A}\|^2}$$

where $\|\cdot\|_{\text{F}}$ is the Frobenius norm and $\|\cdot\|$ is the spectral norm. The stable rank can be viewed as an analytic surrogate for the algebraic rank. Indeed, we may express the two norms in terms of singular values to obtain the relation

$$\text{st.rank}(\mathbf{A}) \leq \text{rank}(\mathbf{A}).$$

In this bound, equality occurs (for example) when the columns of \mathbf{A} are identical or when the columns of \mathbf{A} are orthonormal. As we will see, the stable rank is tightly connected with the number of (strongly) linearly independent columns we can extract from a matrix.

Before we continue, let us instate some regulations. For simplicity we work with real matrices; the complex case requires only minor changes. We say that a matrix is *standardized* when its columns have unit ℓ_2 norm. The j th column of a matrix \mathbf{A} is denoted by \mathbf{a}_j . For a subset τ of column indices, we write \mathbf{A}_τ for the column submatrix indexed by τ . Likewise, given a square matrix \mathbf{H} , the notation $\mathbf{H}_{\tau \times \tau}$ refers to the principal submatrix whose rows and columns are listed in τ . The pseudoinverse \mathbf{D}^\dagger of a diagonal matrix \mathbf{D} is formed by reciprocating the nonzero entries. As usual, we write $\|\cdot\|_p$ for the ℓ_p vector norm. The *condition number* of a matrix is the quantity

$$\kappa(\mathbf{A}) = \max \left\{ \frac{\|\mathbf{Ax}\|_2}{\|\mathbf{Ay}\|_2} : \|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1 \right\}.$$

Finally, upright letters (c, C, K, ...) refer to positive, universal constants that may change from appearance to appearance.

*JAT is with Applied and Computational Mathematics, MC 217-50, California Inst. Technology, Pasadena, CA 91125-5000. E-mail: jtropp@acm.caltech.edu. Supported in part by ONR award no. N00014-08-1-0883.

The first theorem, due to Kashin and Tzafriri, shows that each matrix with standardized columns contains a large column submatrix that has small spectral norm [Ver01, Thm. 2.5].

THEOREM 1.1. (KASHIN–TZAFRIRI) *Suppose \mathbf{A} is standardized. Then there is a set τ of column indices for which*

$$|\tau| \geq \text{st.rank}(\mathbf{A}) \quad \text{and} \quad \|\mathbf{A}_\tau\| \leq C.$$

In fact, much more is true. Combining Theorem 1.1 with the celebrated restricted invertibility result of Bourgain and Tzafriri [BT87, Thm. 1.2], we find that every standardized matrix contains a large column submatrix whose *condition number* is small.

THEOREM 1.2. (BOUGAIN–TZAFRIRI) *Suppose \mathbf{A} is standardized. Then there is a set τ of column indices for which*

$$|\tau| \geq c \cdot \text{st.rank}(\mathbf{A}) \quad \text{and} \quad \kappa(\mathbf{A}_\tau) \leq \sqrt{3}.$$

Theorem 1.2 yields the best general result [BT91, Thm. 1.1] on the Kadison–Singer conjecture, a major open question in operator theory. To display its strength, let us consider two extreme examples.

1. When \mathbf{A} has identical columns, every collection of two or more columns is singular. Theorem 1.2 guarantees a well-conditioned submatrix \mathbf{A}_τ with $|\tau| = 1$, which is optimal.
2. When \mathbf{A} has n orthonormal columns, the full matrix is perfectly conditioned. Theorem 1.2 guarantees a well-conditioned submatrix \mathbf{A}_τ with $|\tau| \geq cn$, which lies within a constant factor of optimal.

Theorem 1.2 uses the stable rank to interpolate between the two extremes. Subsequent research has established that the stable rank is intrinsic to the problem of finding well-conditioned submatrices. We postpone a more detailed discussion of this point until Section 6.

1.1 Contributions. Although Theorems 1.1 and 1.2 would be very useful in computational applications, we cannot regard current proofs as constructive. The goal of this paper is to establish the following novel algorithmic claim.

THEOREM 1.3. *There are randomized, polynomial-time algorithms for producing the sets guaranteed by Theorem 1.1 and by Theorem 1.2.*

This result is significant because no known algorithm for column subset selection is guaranteed to produce a submatrix whose condition number has constant order. See [BDM08] for a recent overview of that literature. The present work has other ramifications with independent interest.

- We develop algorithms for computing the matrix factorizations of Pietsch and Grothendieck, which are regarded as basic instruments in modern functional analysis [Pis86].
- The methods for computing these factorizations lead to approximation algorithms for two NP-hard matrix norms. (See Remarks 3.1 and 5.1.)
- We identify an intriguing connection between Pietsch factorization and the MAXCUT semidefinite program [GW95].

1.2 Overview. We focus on the algorithmic version of the Kashin–Tzafriri theorem because it highlights all the essential concepts while minimizing irrelevant details. Section 2 outlines a proof of this result, emphasizing where new algorithmic machinery is required. The missing link turns out to be a computational method for producing a certain matrix factorization. Section 3 reformulates the factorization problem as an eigenvalue minimization, which can be completed with standard techniques. In Section 4, we exhibit a randomized algorithm that delivers the submatrix promised by Kashin–Tzafriri. In Section 5, we traverse a similar route to develop an algorithmic version of Bourgain–Tzafriri. Section 6 provides more details about the stable rank and describes directions for future work.

2 The Kashin–Tzafriri Theorem.

The proof of the Kashin–Tzafriri theorem proceeds in two steps. First, we select a random set of columns with appropriate cardinality. Second, we use a matrix factorization to identify and remove redundant columns that inflate the spectral norm. The proof gives strong hints about how a computational procedure might work, even though it is not constructive.

2.1 Intuitions. We would like to think that a random submatrix inherits its share of the norm of the entire matrix. In other words, if we were to select a tenth of the columns, we might hope to reduce the norm by a factor of ten. Unfortunately, this intuition is meretricious.

Indeed, random selection does not necessarily reduce the spectral norm at all. The essential reason emerges when we consider the “double identity,” the

$m \times 2m$ matrix $\mathbf{A} = [\mathbf{I} \mid \mathbf{I}]$. Suppose we draw s random columns from \mathbf{A} without replacement. The probability that all s columns are distinct is

$$\begin{aligned} & \frac{2m-2}{2m-1} \times \frac{2m-4}{2m-2} \times \cdots \times \frac{2m-2(s-1)}{2m-(s-1)} \\ & \leq \prod_{j=0}^{s-1} \left(1 - \frac{j}{2m}\right) \\ & \approx \exp \left\{ - \sum_{j=0}^{s-1} \frac{j}{2m} \right\} \approx e^{-s^2/4m}. \end{aligned}$$

Therefore, when $s = \Omega(\sqrt{m})$, sampling almost always produces a submatrix with at least one duplicated column. A duplicated column means that the norm of the submatrix is $\sqrt{2}$, which equals the norm of the full matrix, so no reduction takes place.

Nevertheless, a randomly chosen set of columns from a standardized matrix typically *contains* a large set of columns that has small norm. We will see that the desired subset is exposed by factoring the random submatrix. This factorization, which was invented by Pietsch, is regarded as a basic instrument in modern functional analysis.

2.2 The $(\infty, 2)$ operator norm. Although sampling does not necessarily reduce the spectral norm, it often reduces other matrix norms. Define the natural norm on linear operators from ℓ_∞ to ℓ_2 via the expression

$$\|\mathbf{B}\|_{\infty \rightarrow 2} = \max\{\|\mathbf{B}\mathbf{x}\|_2 : \|\mathbf{x}\|_\infty = 1\}.$$

An immediate consequence is that $\|\mathbf{B}\|_{\infty \rightarrow 2} \leq \sqrt{s} \|\mathbf{B}\|$ for each matrix \mathbf{B} with s columns. Equality can obtain in this bound.

The exact calculation of the $(\infty, 2)$ operator norm is computationally difficult. Results of Rohn [Roh00] imply that there is a class of positive semidefinite matrices for which it is NP-hard to estimate $\|\cdot\|_{\infty \rightarrow 2}$ within an absolute tolerance. Nevertheless, we will see that the norm can be approximated in polynomial time up to a small relative error. (See Remark 3.1.)

As we have intimated, the $(\infty, 2)$ norm can often be reduced by random selection. The following theorem requires some heavy lifting, which we delegate to the technical report [Tro08].

THEOREM 2.1. *Suppose \mathbf{A} is a standardized matrix with n columns. Choose*

$$s \leq \lceil 2 \text{st.rank}(\mathbf{A}) \rceil,$$

and draw a uniformly random subset σ with cardinality s from $\{1, 2, \dots, n\}$. Then

$$\mathbb{E} \|\mathbf{A}_\sigma\|_{\infty \rightarrow 2} \leq 7\sqrt{s}.$$

In particular, $\|\mathbf{A}_\sigma\|_{\infty \rightarrow 2} \leq 8\sqrt{s}$ with probability at least $1/8$.

2.3 Pietsch factorization. We cannot exploit the bound in Theorem 2.1 unless we have a way to connect the $(\infty, 2)$ norm with the spectral norm. To that end, let us recall one of the landmark theorems of functional analysis.

THEOREM 2.2. (PIETSCH FACTORIZATION) *Each matrix \mathbf{B} can be factorized as $\mathbf{B} = \mathbf{T}\mathbf{D}$ where*

- \mathbf{D} is a nonnegative, diagonal matrix with $\text{trace}(\mathbf{D}^2) = 1$, and
- $\|\mathbf{B}\|_{\infty \rightarrow 2} \leq \|\mathbf{T}\| \leq K_P \|\mathbf{B}\|_{\infty \rightarrow 2}$.

This result follows from the little Grothendieck theorem [Pis86, Sec. 5b] and the Pietsch factorization theorem [Pis86, Cor. 1.8]. The standard proof produces the factorization using an abstract separation argument that offers no algorithmic insight. The value of the constant is available.

- When the scalar field is real, we have $K_P(\mathbb{R}) = \sqrt{\pi/2} \approx 1.25$.
- When the scalar field is complex, we have $K_P(\mathbb{C}) = \sqrt{4/\pi} \approx 1.13$.

A major application of Pietsch factorization is to identify a submatrix with controlled spectral norm. The following proposition describes the procedure.

PROPOSITION 2.1. *Suppose \mathbf{B} is a matrix with s columns. Then there is a set τ of column indices for which*

$$|\tau| \geq \frac{s}{2} \quad \text{and} \quad \|\mathbf{B}_\tau\| \leq K_P \sqrt{\frac{2}{s}} \|\mathbf{B}\|_{\infty \rightarrow 2}.$$

Proof. Consider a Pietsch factorization $\mathbf{B} = \mathbf{T}\mathbf{D}$, and define

$$\tau = \{j : d_{jj}^2 \leq 2/s\}.$$

Since $\sum d_{jj}^2 = 1$, Markov's inequality implies that $|\tau| \geq s/2$. We may calculate that

$$\|\mathbf{B}_\tau\| = \|\mathbf{T}\mathbf{D}_\tau\| \leq \|\mathbf{T}\| \cdot \|\mathbf{D}_\tau\| \leq K_P \|\mathbf{B}\|_{\infty \rightarrow 2} \sqrt{2/s}.$$

This completes the proof.

2.4 Proof of Kashin–Tzafriri. With these results at hand, we easily complete the proof of the Kashin–Tzafriri theorem. Suppose \mathbf{A} is a standardized matrix with n columns. We assume that $\text{st.rank}(\mathbf{A}) \leq n/2$. Otherwise, the spectral norm $\|\mathbf{A}\| \leq \sqrt{2}$, so we may select $\tau = \{1, 2, \dots, n\}$.

According to Theorem 2.1, there is a subset σ of column indices for which

$$|\sigma| \geq 2 \text{ st. rank}(\mathbf{A}) \quad \text{and} \quad \|\mathbf{A}_\sigma\|_{\infty \rightarrow 2} \leq 8\sqrt{|\sigma|}.$$

Apply Proposition 2.1 to the matrix $\mathbf{B} = \mathbf{A}_\sigma$ to obtain a subset τ inside σ for which

$$|\tau| \geq \frac{|\sigma|}{2} \quad \text{and} \quad \|\mathbf{B}_\tau\| \leq K_P \sqrt{\frac{2}{|\sigma|}} \|\mathbf{B}\|_{\infty \rightarrow 2}.$$

Since $\mathbf{B}_\tau = \mathbf{A}_\tau$ and $K_P \leq \sqrt{\pi/2}$, these bounds reveal the advertised conclusion:

$$|\tau| \geq \text{st. rank}(\mathbf{A}) \quad \text{and} \quad \|\mathbf{A}_\tau\| < 15.$$

At this point, we take a step back and notice that this proof is nearly algorithmic. It is straightforward to perform the random selection described in Theorem 2.1. Provided that we know a Pietsch factorization of the matrix \mathbf{B} , we can easily carry out the column selection of Proposition 2.1. Therefore, we need only develop an algorithm for computing the Pietsch factorization to reach an effective version of the Kashin–Tzafriri theorem.

3 Pietsch Factorization via Convex Optimization.

The main novelty is to demonstrate that we can produce a Pietsch factorization by solving a convex programming problem. Remarkably, the resulting optimization is the dual of the famous MAXCUT semidefinite program [GW95], for which many polynomial-time algorithms are available.

3.1 Pietsch and eigenvalues. The next theorem, which serves as the basis for our computational method, demonstrates that Pietsch factorizations have an intimate relationship with the eigenvalues of a related matrix. In the sequel, we reserve the letter \mathbf{D} for a nonnegative, diagonal matrix with $\text{trace}(\mathbf{D}^2) = 1$, and we write λ_{\max} for the algebraically maximal eigenvalue of an Hermitian matrix.

THEOREM 3.1. *The factorization $\mathbf{B} = \mathbf{T}\mathbf{D}$ satisfies $\|\mathbf{T}\| \leq \alpha$ if and only if \mathbf{D} satisfies*

$$\lambda_{\max}(\mathbf{B}^*\mathbf{B} - \alpha^2\mathbf{D}^2) \leq 0.$$

In particular, if no \mathbf{D} verifies this bound, then no factorization $\mathbf{B} = \mathbf{T}\mathbf{D}$ admits $\|\mathbf{T}\| \leq \alpha$.

Proof. Assume \mathbf{B} has a factorization $\mathbf{B} = \mathbf{T}\mathbf{D}$ with

$\|\mathbf{T}\| \leq \alpha$. We have the chain of implications

$$\begin{aligned} \mathbf{B} = \mathbf{T}\mathbf{D} &\implies \|\mathbf{B}\mathbf{x}\|_2^2 = \|\mathbf{T}\mathbf{D}\mathbf{x}\|_2^2 && \forall \mathbf{x} \\ &\implies \|\mathbf{B}\mathbf{x}\|_2^2 \leq \alpha^2 \|\mathbf{D}\mathbf{x}\|_2^2 && \forall \mathbf{x} \\ &\implies \mathbf{x}^*\mathbf{B}^*\mathbf{B}\mathbf{x} \leq \alpha^2 \mathbf{x}^*\mathbf{D}^2\mathbf{x} && \forall \mathbf{x} \\ &\implies \mathbf{x}^*(\mathbf{B}^*\mathbf{B} - \alpha^2\mathbf{D}^2)\mathbf{x} \leq 0 && \forall \mathbf{x} \\ &\implies \mathbf{B}^*\mathbf{B} - \alpha^2\mathbf{D}^2 \preceq \mathbf{0}, \end{aligned}$$

where \preceq denotes the semidefinite, or Löwner, ordering on Hermitian matrices.

Conversely, assume we are provided the inequality

$$(3.1) \quad \mathbf{B}^*\mathbf{B} - \alpha^2\mathbf{D}^2 \preceq \mathbf{0}.$$

First, we claim that any zero entry in \mathbf{D} corresponds with a zero column of \mathbf{B} . To check this point, suppose that $d_{jj} = 0$ for an index j . The relation (3.1) requires that

$$0 \geq (\mathbf{B}^*\mathbf{B} - \alpha^2\mathbf{D}^2)_{jj} = \mathbf{b}_j^*\mathbf{b}_j.$$

This inequality is impossible unless $\mathbf{b}_j = \mathbf{0}$. To continue, set $\mathbf{T} = \mathbf{B}\mathbf{D}^\dagger$, and observe that $\mathbf{B} = \mathbf{T}\mathbf{D}$ because the zero entries of \mathbf{D} correspond with zero columns of \mathbf{B} . Therefore, we may factor the diagonal matrix out from (3.1) to reach

$$\mathbf{D}^*(\mathbf{T}^*\mathbf{T} - \alpha^2\mathbf{I})\mathbf{D} \preceq \mathbf{0}.$$

Sylvester's theorem on inertia [HJ85, Thm. 4.5.8] ensures that $\mathbf{T}^*\mathbf{T} - \alpha^2\mathbf{I} \preceq \mathbf{0}$. We conclude that $\|\mathbf{T}\| \leq \alpha$.

3.2 Factorization via optimization. Recall that the maximum eigenvalue is a convex function on the space of Hermitian matrices, so it can be minimized in polynomial time [LO96]. We are led to consider the convex program

$$(3.2) \quad \min \lambda_{\max}(\mathbf{B}^*\mathbf{B} - \alpha^2\mathbf{F}) \quad \text{subject to} \\ \text{trace}(\mathbf{F}) = 1, \mathbf{F} \text{ diagonal, and } \mathbf{F} \geq \mathbf{0}.$$

Owing to Theorem 3.1, there exists a factorization $\mathbf{B} = \mathbf{T}\mathbf{D}$ with $\|\mathbf{T}\| \leq \alpha$ if and only if the value of (3.2) is nonpositive.

Now, if \mathbf{F} is a feasible point of (3.2) with a nonpositive objective value, we can factorize

$$\mathbf{B} = \mathbf{T}\mathbf{D} \quad \text{with} \\ \mathbf{D} = \mathbf{F}^{1/2}, \quad \mathbf{T} = \mathbf{B}\mathbf{D}^\dagger, \quad \text{and} \quad \|\mathbf{T}\| \leq \alpha.$$

In fact, it is not necessary to solve (3.2) to optimality. Suppose \mathbf{B} has s columns, and assume we have identified a feasible point \mathbf{F} with a (positive) objective value η . That is,

$$\lambda_{\max}(\mathbf{B}^*\mathbf{B} - \alpha^2\mathbf{F}) \leq \eta.$$

Rearranging this relation, we reach

$$\lambda_{\max} \left[\mathbf{B}^* \mathbf{B} - (\alpha^2 + \eta s) \tilde{\mathbf{F}} \right] \leq 0 \quad \text{where} \\ \tilde{\mathbf{F}} = \frac{1}{\alpha^2 + \eta s} (\alpha^2 \mathbf{F} + \eta \mathbf{I}).$$

Since $\tilde{\mathbf{F}}$ is positive and diagonal with $\text{trace}(\tilde{\mathbf{F}}) = 1$, we obtain the factorization

$$\mathbf{B} = \mathbf{T} \mathbf{D} \quad \text{with} \quad \mathbf{D} = \tilde{\mathbf{F}}^{1/2}, \\ \mathbf{T} = \mathbf{B} \mathbf{D}^{-1}, \quad \text{and} \quad \|\mathbf{T}\| \leq \sqrt{\alpha^2 + \eta s}.$$

To select a target value for the parameter α , we look to the proof of the Kashin–Tzafriri theorem. If \mathbf{B} has s columns, then $\alpha = 8K_P \sqrt{s}$ is an appropriate choice. Furthermore, since the argument only uses the bound $\|\mathbf{T}\| = O(\sqrt{s})$, it suffices to solve (3.2) with precision $\eta = O(1)$.

3.3 Other formulations. In a general setting, a target value for α is not likely to be available. Let us exhibit an alternative formulation of (3.2) that avoids this inconvenience:

$$(3.3) \quad \min \lambda_{\max}(\mathbf{B}^* \mathbf{B} - \mathbf{E}) + \text{trace}(\mathbf{E}) \\ \text{subject to} \quad \mathbf{E} \text{ diagonal, } \mathbf{E} \geq \mathbf{0}.$$

Suppose α_* is the minimal value of $\|\mathbf{T}\|$ achievable in any Pietsch factorization $\mathbf{B} = \mathbf{T} \mathbf{D}$. It can be shown that α_*^2 is the value of (3.3) and that each optimizer \mathbf{E}_* satisfies $\text{trace}(\mathbf{E}_*) = \alpha_*^2$. As such, we can construct an optimal Pietsch factorization from a minimizer:

$$\mathbf{B} = \mathbf{T} \mathbf{D} \quad \text{with} \quad \mathbf{D} = (\mathbf{E}_* / \text{trace}(\mathbf{E}_*))^{1/2}, \\ \mathbf{T} = \mathbf{B} \mathbf{D}^\dagger, \quad \text{and} \quad \|\mathbf{T}\| = \alpha_*.$$

The dual of (3.3) is the semidefinite program

$$(3.4) \quad \max \langle \mathbf{B}^* \mathbf{B}, \mathbf{Z} \rangle \\ \text{subject to} \quad \text{diag}(\mathbf{Z}) = \mathbf{I} \text{ and } \mathbf{Z} \succcurlyeq \mathbf{0}.$$

This is the famous MAXCUT semidefinite program [GW95]. We find an unexpected connection between Pietsch factorization and the problem of partitioning nodes of a graph.

Given a dual optimum, we can easily construct a primal optimum by means of the complementary slackness condition [Ali95, Thm. 2.10]. Indeed, each feasible optimal pair $(\mathbf{E}_*, \mathbf{Z}_*)$ satisfies $\mathbf{Z}_*(\mathbf{B}^* \mathbf{B} - \mathbf{E}_*) = \mathbf{0}$. Examining the diagonal elements of this matrix equation, we find that

$$\mathbf{E}_* = \text{diag}(\mathbf{E}_*) = \text{diag}(\mathbf{Z} \mathbf{E}_*) = \text{diag}(\mathbf{Z}_* \mathbf{B}^* \mathbf{B})$$

owing to the constraint $\text{diag}(\mathbf{Z}_*) = \mathbf{I}$. Obtaining a dual optimum from a primal optimum, however, requires more ingenuity.

REMARK 3.1. According to Theorem 2.2 and the discussion here, the optimal value of (3.3) overestimates $\|\mathbf{B}\|_{\infty \rightarrow 2}^2$ by a multiplicative factor no greater than K_P^2 . As a result, the optimization problem (3.3) can be used to design an approximation algorithm for $(\infty, 2)$ norms.

3.4 Algorithmic aspects. The purpose of this paper is not to rehash methods for solving a standard optimization problem, so we keep this discussion brief. It is easy to see that (3.2) can be framed as a (nonsmooth) convex optimization over the probability simplex. The technical report [Tro08] outlines an elegant technique, called Entropic Mirror Descent [BT03], designed specifically for this class of problems. Although the EMD algorithm is (theoretically) not the most efficient approach to (3.2), preliminary experiments suggest that its empirical performance rivals more sophisticated techniques.

For a concrete time bound, we refer to Alizadeh’s work on primal–dual potential reduction methods for semidefinite programming [Ali95]. When \mathbf{B} has dimension $m \times s$, the cost of forming $\mathbf{B}^* \mathbf{B}$ is at most $O(s^2 m)$. Then the cost of solving (3.4) is no more than $\tilde{O}(s^{3.5})$, where the tilde indicates that log-like factors are suppressed.

4 An Algorithm for Kashin–Tzafriri.

At this point, we have amassed the matériel necessary to deploy an algorithm that constructs the set τ promised by the Kashin–Tzafriri theorem. The procedure appears as Algorithm 1. The following result describes its performance.

THEOREM 4.1. Suppose \mathbf{A} is an $m \times n$ standardized matrix. With probability at least $4/5$, Algorithm 1 produces a set $\tau = \tau_*$ of column indices for which

$$|\tau| \geq \text{st.rank}(\mathbf{A}) \quad \text{and} \quad \|\mathbf{A}_\tau\| \leq 15.$$

The computational cost is bounded by $\tilde{O}(|\tau|^2 m + |\tau|^{3.5})$.

Remarkably, Algorithm 1 is sublinear in the size of the matrix when $\text{st.rank}(\mathbf{A}) = o(n^{1/3.5})$. Better methods for solving (3.2) would strengthen this bound.

Proof. According to Section 2, the procedure NORM-REDUCE has failure probability less than $7/8$ when $s \leq 2 \text{st.rank}(\mathbf{A})$. The probability the inner loop

fails to produce an acceptable set τ_* of size $s/2$ is at most $(7/8)^{8 \log_2(s)}$. So the probability the algorithm fails before $s > 2 \text{st.rank}(\mathbf{A})$ is at most

$$\sum_{j=2}^{\infty} (7/8)^{8j} = \frac{(7/8)^{16}}{1 - (7/8)^8} < 0.2.$$

With constant probability, we obtain a set τ_* with cardinality at least $\text{st.rank}(\mathbf{A})$.

The cost of the procedure NORM-REDUCE is dominated by the cost of the Pietsch factorization, which is $\tilde{O}(s^2 m + s^{3.5})$ for a fixed s . Summing over s and k , we find that the total cost of all the invocations of NORM-REDUCE is dominated (up to logarithmic factors) by the cost of the final invocation, during which the parameter $s \leq 2|\tau_*|$.

An estimate of the spectral norm of \mathbf{A}_τ can be obtained as a by-product of solving (3.2). Indeed, Proposition 2.1 and the discussion in Section 3.2 show that we can bound the spectral norm in terms of the parameter α and the objective value obtained in (3.2).

Algorithm 1: Constructive version of Kashin–Tzafriri theorem

KT(\mathbf{A})
Input: Standardized matrix \mathbf{A} with n columns
Output: A subset τ_* of $\{1, 2, \dots, n\}$
Description: Produces τ_* such that $|\tau_*| \geq \text{st.rank}(\mathbf{A})$ and $\|\mathbf{A}_\tau\| \leq 15$ w.p. $4/5$

```

1   $\tau_* = \{1\}$ 
2  for  $s = 4, 8, 16, \dots, n$ 
3      for  $k = 1, 2, 3, \dots, 8 \log_2 s$ 
4           $\tau = \text{NORM-REDUCE}(\mathbf{A}, s)$ 
5          if  $\|\mathbf{A}_\tau\| \leq 15$  then  $\tau_* = \tau$  and break
6  if  $|\tau_*| < s$  then exit

```

NORM-REDUCE(\mathbf{A}, s)
Input: Standardized matrix \mathbf{A} with n columns, a positive integer s
Output: A subset τ of $\{1, 2, \dots, n\}$

```

1  Draw a uniformly random set  $\sigma$  with cardinality  $s$  from  $\{1, 2, \dots, n\}$ 
2  Solve (3.2) with  $\mathbf{B} = \mathbf{A}_\sigma$  and  $\alpha = 8K_P\sqrt{s}$  to obtain a factorization  $\mathbf{B} = \mathbf{T}\mathbf{D}$ 
3  Return  $\tau = \{j \in \sigma : d_{jj}^2 \leq 2/s\}$ 

```

5 The Bourgain–Tzafriri Theorem.

Our proof of the Bourgain–Tzafriri theorem is almost identical in structure with the proof of the Kashin–Tzafriri theorem. This streamlined argument ap-

pears to be simpler than all previously published approaches, but it contains no significant conceptual innovations. Our discussion culminates in an algorithm remarkably similar to Algorithm 1.

5.1 Preliminary results. Suppose \mathbf{A} is a standardized matrix with n columns. We will work instead with a related matrix $\mathbf{H} = \mathbf{A}^* \mathbf{A} - \mathbf{I}$, which is called the *hollow Gram matrix*. The advantage of considering the hollow Gram matrix is that we can perform column selection on \mathbf{A} simply by reducing the norm of \mathbf{H} .

PROPOSITION 5.1. *Suppose \mathbf{A} is a standardized matrix with hollow Gram matrix \mathbf{H} . If τ is a set of column indices for which $\|\mathbf{H}_{\tau \times \tau}\| \leq 0.5$, then $\kappa(\mathbf{A}_\tau) \leq \sqrt{3}$.*

Proof. The hypothesis $\|\mathbf{H}_{\tau \times \tau}\| \leq 0.5$ implies that the eigenvalues of $\mathbf{H}_{\tau \times \tau}$ lie in the range $[-0.5, 0.5]$. Since $\mathbf{H}_{\tau \times \tau} = \mathbf{A}_\tau^* \mathbf{A}_\tau - \mathbf{I}$, the eigenvalues of $\mathbf{A}_\tau^* \mathbf{A}_\tau$ fall in the interval $[0.5, 1.5]$. An equivalent condition is that $0.5 \leq \|\mathbf{A}_\tau \mathbf{x}\|_2^2 \leq 1.5$ whenever $\|\mathbf{x}\|_2 = 1$. We conclude that

$$\begin{aligned} \kappa(\mathbf{A}_\tau) &= \max \left\{ \frac{\|\mathbf{A}_\tau \mathbf{x}\|_2}{\|\mathbf{A}_\tau \mathbf{y}\|_2} : \|\mathbf{x}\|_2 = \|\mathbf{y}\|_2 = 1 \right\} \\ &\leq \sqrt{\frac{1.5}{0.5}} = \sqrt{3}. \end{aligned}$$

Thus, a norm bound for $\mathbf{H}_{\tau \times \tau}$ yields a condition number bound for \mathbf{A}_τ .

As we mentioned before, random selection may reduce other norms even if it does not reduce the spectral norm. Define the natural norm on linear maps from ℓ_∞ to ℓ_1 by the formula

$$\|\mathbf{G}\|_{\infty \rightarrow 1} = \max\{\|\mathbf{G}\mathbf{x}\|_1 : \|\mathbf{x}\|_\infty = 1\}.$$

This norm is closely related to the cut norm, which plays a starring role in graph theory [AN04]. For a general $s \times s$ matrix \mathbf{G} , the best inequality between the $(\infty, 1)$ norm and the spectral norm is $\|\mathbf{G}\|_{\infty \rightarrow 1} \leq s \|\mathbf{G}\|$. Rohn [Roh00] has established that there is a class of positive semidefinite, integer matrices for which it is NP-hard to determine the $(\infty, 1)$ norm within an absolute tolerance of $1/2$. Nevertheless, it can be approximated within a small relative factor in polynomial time [AN04].

The $(\infty, 1)$ norm decreases when we randomly sample a principal submatrix. The following result, established in [Tro08] is a direct consequence of Rudelson and Vershynin’s work on the cut norm of random submatrices [RV07, Thm. 1.5].

THEOREM 5.1. Suppose \mathbf{A} is an n -column standardized matrix with hollow Gram matrix \mathbf{H} . Choose

$$s \leq \lceil c \cdot \text{st.rank}(\mathbf{A}) \rceil,$$

and draw a uniformly random subset σ with cardinality s from $\{1, 2, \dots, n\}$. Then

$$\mathbb{E} \|\mathbf{H}_{\sigma \times \sigma}\|_{\infty \rightarrow 1} \leq \frac{s}{9}.$$

In particular, $\|\mathbf{H}_{\sigma \times \sigma}\|_{\infty \rightarrow 1} \leq s/8$ with probability at least $1/9$.

To connect the $(\infty, 1)$ norm with the spectral norm, we call on the celebrated factorization of Grothendieck [Pis86, p. 56].

THEOREM 5.2. (GROTHENDIECK FACTORIZATION) Each matrix \mathbf{G} can be factorized as $\mathbf{G} = \mathbf{D}_1 \mathbf{T} \mathbf{D}_2$ where

1. \mathbf{D}_i is a nonnegative, diagonal matrix with $\text{trace}(\mathbf{D}_i^2) = 1$ for $i = 1, 2$, and
2. $\|\mathbf{G}\|_{\infty \rightarrow 1} \leq \|\mathbf{T}\| \leq K_G \|\mathbf{G}\|_{\infty \rightarrow 1}$.

When \mathbf{G} is Hermitian, we may take $\mathbf{D}_1 = \mathbf{D}_2$.

The precise value of the Grothendieck constant K_G remains an outstanding open question, but it is known to depend on the scalar field [Pis86, Sec. 5e].

- When the scalar field is real, $1.57 \leq \pi/2 \leq K_G(\mathbb{R}) \leq \pi/(2 \log(1 + \sqrt{2})) \leq 1.79$.
- When the scalar field is complex, $1.33 \leq K_G(\mathbb{C}) \leq 1.41$.

For positive semidefinite \mathbf{G} , the real (resp., complex) Grothendieck constant equals the square of the real (resp., complex) Pietsch constant because $\|\mathbf{B}^* \mathbf{B}\|_{\infty \rightarrow 1} = \|\mathbf{B}\|_{\infty \rightarrow 2}^2$.

The following proposition describes the role of the Grothendieck factorization in the selection of submatrices with controlled spectral norm.

PROPOSITION 5.2. Suppose \mathbf{G} is an $s \times s$ Hermitian matrix. There is a set τ of column indices for which

$$|\tau| \geq \frac{s}{2} \quad \text{and} \quad \|\mathbf{G}_{\tau \times \tau}\| \leq \frac{2K_G}{s} \|\mathbf{G}\|_{\infty \rightarrow 1}.$$

Proof. Consider a Grothendieck factorization $\mathbf{G} = \mathbf{D} \mathbf{T} \mathbf{D}$, and identify $\tau = \{j : d_{jj}^2 \leq s/2\}$. The remaining details echo the proof of Proposition 2.1.

5.2 Proof of Bourgain–Tzafriri. Suppose \mathbf{A} is a standardized matrix with n columns, and consider its hollow Gram matrix \mathbf{H} . Theorem 5.1 provides a set σ for which

$$|\sigma| \geq c \cdot \text{st.rank}(\mathbf{A}) \quad \text{and} \quad \|\mathbf{H}_{\sigma \times \sigma}\|_{\infty \rightarrow 1} \leq \frac{s}{8}.$$

Apply Proposition 5.2 to the $s \times s$ matrix $\mathbf{G} = \mathbf{H}_{\sigma \times \sigma}$ to obtain a further subset τ inside σ with

$$|\tau| \geq \frac{s}{2} \quad \text{and} \quad \|\mathbf{G}_{\tau \times \tau}\| \leq \frac{2K_G}{s} \|\mathbf{G}\|_{\infty \rightarrow 1}.$$

Since $2K_G < 4$ and $\mathbf{H}_{\tau \times \tau} = \mathbf{G}_{\tau \times \tau}$, we determine that

$$|\tau| \geq \frac{c}{2} \cdot \text{st.rank}(\mathbf{A}) \quad \text{and} \quad \|\mathbf{H}_{\tau \times \tau}\| \leq 0.5.$$

In view of Proposition 5.1, we conclude $\kappa(\mathbf{A}_\tau) \leq \sqrt{3}$.

Now, take another step back and notice that this argument is nearly algorithmic. The random selection of σ can easily be implemented in practice, even though the proof does not specify the value of c . Given a Grothendieck factorization $\mathbf{G} = \mathbf{D} \mathbf{T} \mathbf{D}$, it is straightforward to identify the subset τ . The challenge, as before, is to produce the factorization.

5.3 Grothendieck factorization via convex optimization. As with the Pietsch factorization, the Grothendieck factorization can be identified from the solution to a convex program.

THEOREM 5.3. Suppose \mathbf{G} is Hermitian. The factorization $\mathbf{G} = \mathbf{D} \mathbf{T} \mathbf{D}$ satisfies $\|\mathbf{T}\| \leq \alpha$ if and only if \mathbf{D} satisfies

$$(5.5) \quad \lambda_{\max} \begin{bmatrix} -\alpha \mathbf{D}^2 & \mathbf{G} \\ \mathbf{G} & -\alpha \mathbf{D}^2 \end{bmatrix} \leq 0.$$

In particular, if no \mathbf{D} verifies this bound, then no factorization $\mathbf{G} = \mathbf{D} \mathbf{T} \mathbf{D}$ admits $\|\mathbf{T}\| \leq \alpha$.

Proof. To check the forward implication, we essentially repeat the argument we used in Theorem 3.1 for the Pietsch case. This reasoning yields the pair of relations

$$\mathbf{G} - \alpha \mathbf{D}^2 \preceq \mathbf{0} \quad \text{and} \quad -\mathbf{G} - \alpha \mathbf{D}^2 \preceq \mathbf{0}.$$

Together, these two relations are equivalent with (5.5) because

$$\begin{bmatrix} -\alpha \mathbf{D}^2 & \mathbf{G} \\ \mathbf{G} & -\alpha \mathbf{D}^2 \end{bmatrix} = \frac{1}{2} \begin{bmatrix} \mathbf{I} & \mathbf{I} \\ -\mathbf{I} & \mathbf{I} \end{bmatrix}^* \begin{bmatrix} \mathbf{G} - \alpha \mathbf{D}^2 & \\ & -\mathbf{G} - \alpha \mathbf{D}^2 \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{I} \\ -\mathbf{I} & \mathbf{I} \end{bmatrix}.$$

To prove the reverse implication, we assume that (5.5) holds. First, we must check that $d_{jj} = 0$ implies that $\mathbf{g}_j = \mathbf{0}$. To verify this claim, observe that

$$\begin{aligned} 0 &\geq \begin{bmatrix} \alpha \\ \mathbf{g}_j \end{bmatrix}^* \begin{bmatrix} 0 & \mathbf{g}_j^* \\ \mathbf{g}_j & -\alpha \mathbf{D}^2 \end{bmatrix} \begin{bmatrix} \alpha \\ \mathbf{g}_j \end{bmatrix} \\ &= \alpha \left(2 \|\mathbf{g}_j\|_2^2 - \mathbf{g}_j^* \mathbf{D}^2 \mathbf{g}_j \right) \geq \alpha \|\mathbf{g}_j\|_2^2 \end{aligned}$$

because $\text{trace}(\mathbf{D}^2) = 1$. Therefore, we may construct a Grothendieck factorization $\mathbf{G} = \mathbf{D}\mathbf{T}\mathbf{D}$ with $\|\mathbf{T}\| \leq \alpha$ by setting $\mathbf{T} = \mathbf{D}^\dagger \mathbf{G} \mathbf{D}^\dagger$.

This discussion leads us to frame the eigenvalue minimization problem

$$(5.6) \quad \min \lambda_{\max} \begin{bmatrix} -\alpha \mathbf{F} & \mathbf{G} \\ \mathbf{G} & -\alpha \mathbf{F} \end{bmatrix} \quad \text{subject to} \\ \text{trace}(\mathbf{F}) = 1, \mathbf{F} \text{ diagonal}, \mathbf{F} \geq \mathbf{0}.$$

Owing to Theorem 5.3, there is a factorization $\mathbf{G} = \mathbf{D}\mathbf{T}\mathbf{D}$ with $\|\mathbf{T}\| \leq \alpha$ if and only if the value of (5.6) is nonpositive.

As in Section 3.2, we can easily construct Grothendieck factorizations from (imprecise) solutions to the problem (5.6). The proof of Bourgain–Tzafriri suggests that an appropriate value for the parameter $\alpha = s/4$. Furthermore, we do not need to solve (5.6) to optimality to obtain the required information. Indeed, it suffices to produce a feasible point with an objective value of $O(1)$.

To solve (5.6) in practice, we again propose the Entropic Mirror Descent algorithm [BT03]. To provide a concrete bound on the computational cost, we remark that, when \mathbf{A}_τ has dimension $m \times s$, forming $\mathbf{G} = \mathbf{A}_\tau^* \mathbf{A}_\tau - \mathbf{I}$ costs at most $O(s^2 m)$, and Alizadeh’s interior-point method [Ali95] requires $\tilde{O}(s^{3.5})$ time.

REMARK 5.1. For symmetric \mathbf{G} , Theorem 5.2 shows that the norm $\|\mathbf{G}\|_{\infty \rightarrow 1}$ is approximated within a factor K_G by the least α for which (5.6) has a nonpositive value. A natural reformulation of (5.6) can identify this value of α automatically (cf. Section 3.3). For nonsymmetric \mathbf{G} , similar optimization problems arise. These ideas yield approximation algorithms for the $(\infty, 1)$ norm.

5.4 An algorithm for Bourgain–Tzafriri. We are prepared to state our algorithm for producing the set τ described by the Bourgain–Tzafriri theorem. The procedure appears as Algorithm 2. Note the striking similarity with Algorithm 1. The following result describes the performance of the algorithm. We omit the proof, which parallels that of Theorem 4.1.

THEOREM 5.4. Suppose \mathbf{A} is an $m \times n$ standardized matrix. With probability at least $3/4$, Algorithm 2 produces a set $\tau = \tau_*$ of column indices for which

$$|\tau| \geq c \cdot \text{st.rank}(\mathbf{A}) \quad \text{and} \quad \kappa(\mathbf{A}_\tau) \leq \sqrt{3}.$$

The computational cost is bounded by $\tilde{O}(|\tau|^2 m + |\tau|^{3.5})$.

Algorithm 2: Constructive version of Bourgain–Tzafriri Theorem

BT(\mathbf{A})
Input: Standardized matrix \mathbf{A} with n columns
Output: A subset τ_* of $\{1, 2, \dots, n\}$
Description: Produces τ_* such that $|\tau_*| \geq \text{st.rank}(\mathbf{A})/2$ and $\kappa(\mathbf{A}_{\tau_*}) \leq \sqrt{3}$ w.p. $3/4$

```

1   $\tau_* = \{1\}$ 
2  for  $s = 4, 8, 16, \dots, n$ 
3      for  $k = 1, 2, 3, \dots, 8 \log_2 s$ 
4           $\tau = \text{COND-REDUCE}(\mathbf{A}, s)$ 
5          if  $\kappa(\mathbf{A}_\tau) \leq \sqrt{3}$  then  $\tau_* = \tau$  and break
6  if  $|\tau_*| < s$  then exit
```

COND-REDUCE(\mathbf{A}, s)
Input: Standardized matrix \mathbf{A} with n columns, a parameter s
Output: A subset τ of $\{1, 2, \dots, n\}$

```

1  Draw a uniformly random set  $\sigma$  with cardinality  $s$  from  $\{1, 2, \dots, n\}$ 
2  Solve (5.6) with  $\mathbf{G} = \mathbf{A}_\sigma^* \mathbf{A}_\sigma - \mathbf{I}$  and  $\alpha = s/4$  to obtain factorization  $\mathbf{G} = \mathbf{D}\mathbf{T}\mathbf{D}$ 
3  Return  $\tau = \{j \in \sigma : d_{jj}^2 \leq 2/s\}$ 
```

6 Future Directions.

After the initial work [BT87], additional research has clarified the role of the stable rank. We highlight a positive result of Vershynin [Ver01, Cor. 7.1] and a negative result of Szarek [Sza90, Thm. 1.2] which together imply that the stable rank describes *precisely* how large a well-conditioned column submatrix can in general exist. See [Ver01, Sec. 5] for a more detailed discussion.

THEOREM 6.1. (VERSHYNIN 2001) Fix $\varepsilon > 0$. For each matrix \mathbf{A} , there is a set τ of column indices for which

$$|\tau| \geq (1 - \varepsilon) \cdot \text{st.rank}(\mathbf{A}) \quad \text{and} \quad \kappa(\mathbf{A}_\tau) \leq C(\varepsilon).$$

THEOREM 6.2. (SZAREK) There is a sequence $\{\mathbf{A}(n)\}$ of matrices of increasing dimension for

which

$$|\tau| = \text{st.rank}(\mathbf{A}) \implies \kappa(\mathbf{A}_\tau) = \omega(1).$$

Vershynin's proof constructs the set τ in Theorem 6.1 with a complicated iteration that interleaves the Kashin–Tzafriri theorem and the Bourgain–Tzafriri theorem. We believe that the argument can be simplified substantially and developed into a column selection algorithm. This achievement might lead to a new method for performing rank-revealing factorizations, which could have a significant impact on the practice of numerical linear algebra.

Acknowledgments.

The author thanks Ben Recht for valuable discussions about eigenvalue minimization.

References

- [Ali95] F. Alizadeh. Interior-point methods in semidefinite programming with applications to combinatorial optimization. *SIAM J. Optimization*, 5(1):13–51, Feb. 1995.
- [AN04] N. Alon and A. Naor. Approximating the cut norm via Grothendieck's inequality. In *Proc. 36th Ann. ACM Symposium on Theory of Computing (STOC)*, pages 72–80, Chicago, 2004.
- [BDM08] C. Boutsidis, P. Drineas, and M. Mahoney. On selecting exactly k columns from a matrix. Submitted for publication, 2008.
- [BT87] J. Bourgain and L. Tzafriri. Invertibility of “large” submatrices with applications to the geometry of Banach spaces and harmonic analysis. *Israel J. Math.*, 57(2):137–224, 1987.
- [BT91] J. Bourgain and L. Tzafriri. On a problem of Kadison and Singer. *J. reine angew. Math.*, 420:1–43, 1991.
- [BT03] A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. *Operations Res. Lett.*, 31:167–175, 2003.
- [GE96] M. Gu and S. Eisenstat. Efficient algorithms for computing a strong rank-revealing QR factorization. *SIAM J. Sci. Comput.*, 17(4):848–869, Jul. 1996.
- [GW95] M. X. Goemans and D. P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *J. Assoc. Comput. Mach.*, 42:1115–1145, 1995.
- [HJ85] R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge Univ. Press, 1985.
- [LO96] A. S. Lewis and M. L. Overton. Eigenvalue optimization. *Acta Numerica*, 5:149–190, 1996.
- [Pis86] G. Pisier. *Factorization of linear operators and geometry of Banach spaces*. Number 60 in CBMS Regional Conference Series in Mathematics. AMS, Providence, 1986. Reprinted with corrections, 1987.
- [Roh00] J. Rohn. Computing the norm $\|A\|_{\infty,1}$ is NP-hard. *Linear and Multilinear Algebra*, 47:195–204, 2000.
- [RV07] M. Rudelson and R. Vershynin. Sampling from large matrices: An approach through geometric functional analysis. *J. Amer. Comput. Soc.*, 54(4):Article 21, pp. 1–19, Jul. 2007.
- [Sza90] S. Szarek. Spaces with large distance from ℓ_∞^n and random matrices. *Amer. J. Math.*, 112(6):899–942, Dec. 1990.
- [Tro08] J. A. Tropp. Column subset selection, matrix factorization, and eigenvalue optimization. ACM Technical Report 2008-02, California Institute of Technology, 2008.
- [Ver01] R. Vershynin. Johns decompositions: Selecting a large part. *Israel J. Math.*, 122:253–277, 2001.